# Cube Pruning as Heuristic Search
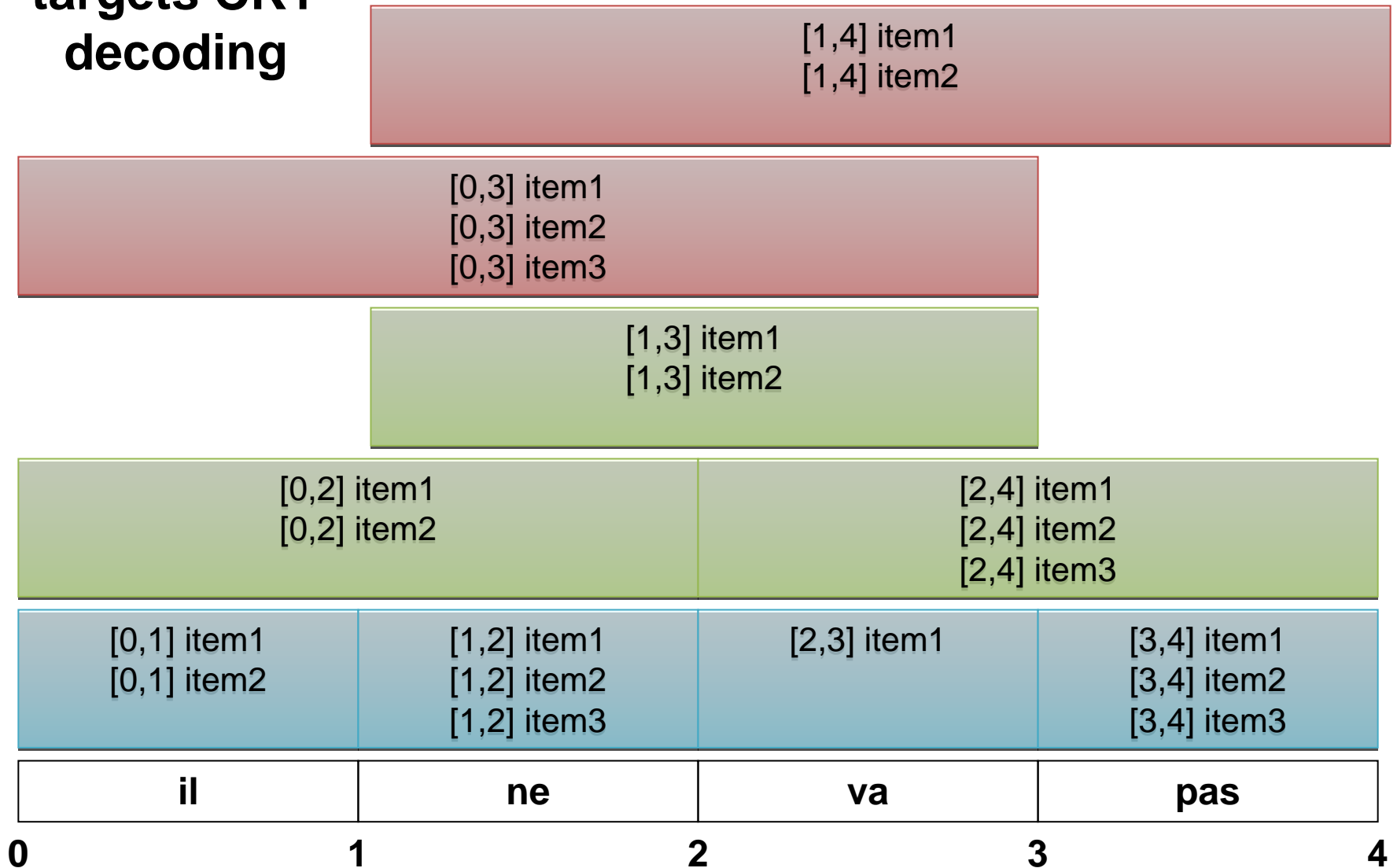
Mark Hopkins and Greg Langmead
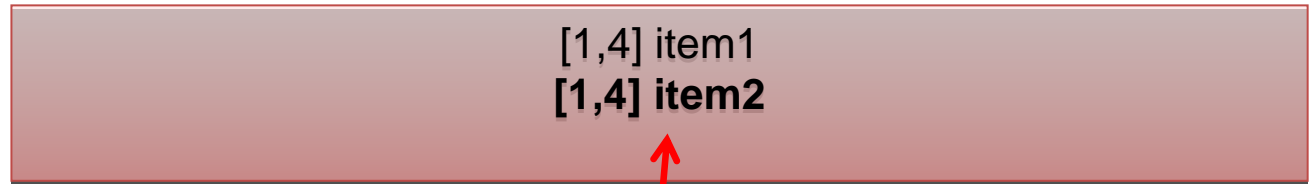
Language Weaver, Inc.

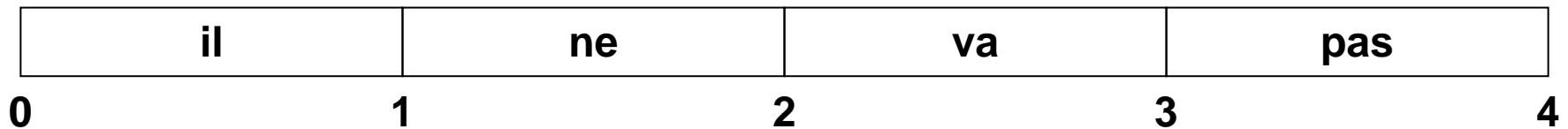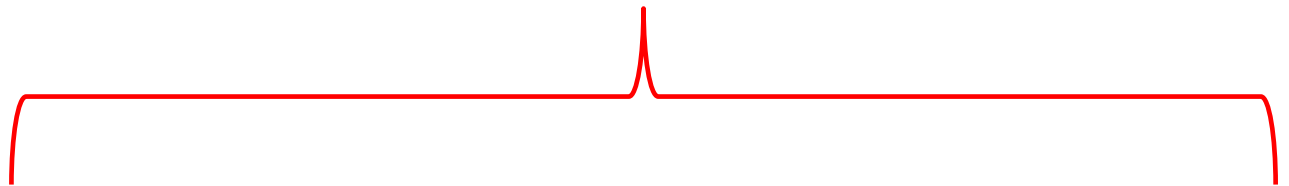# Motivation: Speed vs Accuracy

- Accuracy (e.g., BLEU) is very important
- But time is money
  - for the customer: throughput
  - for LW: cpu time on servers for SaaS solution
- Our customers expect 1000-3000 words per minute in one thread
  - and linear speedup with multiple threads
- That's 0.3-1.0 seconds per sentence
- Can syntax be viable at product speeds?

LANGUAGE WEAVER

# Cube pruning targets CKY decoding

| [1,4] item1 |
| :---: |
| [1,4] item2 |

| [0,3] item1 |
| :---: |
| [0,3] item2 |
| [0,3] item3 |

| [1,3] item1 |
| :---: |
| [1,3] item2 |

| [0,2] item1 [0,2] item2 | [2,4] item1 [2,4] item2 [2,4] item3 |
| :---: | :---: |

| [0,1] item1 [0,1] item2 | [1,2] item1 [1,2] item2 [1,2] item3 | [2,3] item1 | [3,4] item1 [3,4] item2 [3,4] item3 |
| :---: | :---: | :---: | :---: |
| **il** | **ne** | **va** | **pas** |

0    1    2    3    4

LANGUAGE WEAVER

[1,4] item1
**[1,4] item2**

**an item encodes a (set of) translation(s) of
a partial sentence**

| il | ne | va | pas |
|---|---|---|---|

**0**               **1**               **2**               **3**               **4**

LANGUAGE
WEAVER

# How are items created?

????

| | |
|---|---|
| [2,4] item1<br>[2,4] item2<br>[2,4] item3<br>[2,4] item4 | [4,5] item1<br>[4,5] item2<br>[4,5] item3 |

| | |
|---|---|
| [2,3] item1<br>[2,3] item2<br>[2,3] item3 | [3,5] item1<br>[3,5] item2<br>[3,5] item3<br>[3,5] item4 |

**2**       **3**       **4**       **5**

LANGUAGE WEAVER

# Items are created by combining items from complementary spans.

| | | | |
|---|---|---|---|
| **[2,5] item1** | | | |

| | |
|---|---|
| [2,4] item1<br>[2,4] item2<br>**[2,4] item3**<br>[2,4] item4 | [4,5] item1<br>**[4,5] item2**<br>[4,5] item3 |

| | |
|---|---|
| [2,3] item1<br>[2,3] item2<br>[2,3] item3 | [3,5] item1<br>[3,5] item2<br>[3,5] item3<br>[3,5] item4 |

**2**          **3**          **4**          **5**

LANGUAGE WEAVER

# What is an item?

## [ 2,4 , NP , the*car ]

**an item consists of three parts:**

span     postcondition     carry

LANGUAGE WEAVER

# What is an item?

## [ 2,4 , **NP** , the*car ]

**an item consists of three parts:**

**span**     **postcondition**     **carry**

LANGUAGE WEAVER

# What is an item?

## [ 2,4 , NP , the*car ]

**an item consists of three parts:**

span     postcondition     carry

LANGUAGE
WEAVER

# CKY Item Creation

postcondition     carry

postcondition

$[2,3,A,a*b]$

$[3,5,B,b*a]$

$B \rightarrow < A_0 \; B_1 \; , \; B_1 \; A_0 >$

preconditions

$[2,5,B,b*b]$

$$\left[ \begin{array}{c} \text{cost( subitem1)} \\ \text{cost( subitem2)} \\ \text{cost( rule )} \\ + \text{interaction( rule, subitems)} \\ \hline \text{cost( new item )} \end{array} \right]$$

# The Item Creation Problem

**100010001000000001000**

**…then item creation can still take a very long time**

**even if we just store 1000 items here**

**…and 1000 items here**

**…and we have only 1000s of grammar rules**

????

[2,4] item1
[2,4] item2
[2,4] item3
[2,4] item4

[4,5] item1
[4,5] item2
[4,5] item3

[2,3] item1
[2,3] item2
[2,3] item3

[3,5] item1
[3,5] item2
[3,5] item3
[3,5] item4

2    3    4    5

LANGUAGE WEAVER

# The Item Creation Problem

**is there a better way to enumerate the 1000 items of lowest cost for this span without going through the millions of candidate items and taking the best 1000?**

**this is the problem that cube pruning addresses**

| ???? | |
|------|--|
| 1000 | 1000 |

| 1000 | 1000 |
|------|------|

2                         3                        4                        5

LANGUAGE WEAVER

# A demonstration of incremental CKY item creation for span [2,5]

**We want:**

[2,3,A,a*b]

[3,5,B,b*a]

B → < $A_0$ $B_1$ , $B_1$ $A_0$ >

_____

[2,5,B,b*b]

**So far we have:**

[?,?,?,?]

[?,?,?,?]

? → < $?_0$ $?_1$ , ? ? >

_____

[?,?,?,?]

[2,4] [4,5]  [2,3] [3,5]

**We want:**

[2,3,A,a*b]

[3,5,B,b*a]

B → < A₀ B₁ , B₁ A₀ >

———————————

[2,5,B,b*b]

**So far we have:**

[2,3,?,?]

[3,5,?,?]

? → < ?₀ ?₁ , ? ? >

———————————

[2,5,?,?]

LANGUAGE WEAVER

[2,4] [4,5]          [2,3] [3,5]

A          B

**We want:**

[2,3,A,a*b]

[3,5,B,b*a]

B → < A$_0$ B$_1$ , B$_1$ A$_0$ >

_____

[2,5,B,b*b]

**So far we have:**

[2,3,A,?]

[3,5,?,?]

? → < A$_0$ ?$_1$ , ? ? >

_____

[2,5,?,?]

[2,4] [4,5]          [2,3] [3,5]

A          B

A          B

**We want:**

[2,3,A,a*b]

[3,5,B,b*a]

B → < A₀ B₁ , B₁ A₀ >

_____

[2,5,B,b*b]

**So far we have?**

[2,3,A,?]

[3,5,B,?]

? → < A₀ B₁ , ? ? >

_____

[2,5,?,?]

LANGUAGE WEAVER

[2,4] [4,5]   [2,3] [3,5]

A       B

A       B

accept rule(A,B,1)?

y       n

**rule(A,B,k) is the kth lowest cost rule whose preconditions are <A,B>**

**We want:**

[2,3,A,a*b]

[3,5,B,b*a]

B → < A₀ B₁ , B₁ A₀ >

———————————

[2,5,B,b*b]

**So far we have?**

[2,3,A,?]

[3,5,B,?]

B → < A₀ B₁ , B₁ A₀ >

———————————

[2,5,B,?]

LANGUAGE WEAVER

**[2,4] [4,5]**     **[2,3] [3,5]**

A     B

A     B

accept rule(A,B,1)?

y     n

accept item(2,3,A,1)?

y     n

**item(2,3,A,k) is the kth lowest cost item of span [2,3] whose postcondition is A**

## We want:

**[2,3,A,a*b]**

**[3,5,B,b*a]**

**B → < A_0 B_1 , B_1 A_0 >**

---

**[2,5,B,b*b]**

## So far we have:

**[2,3,A,a*b]**

**[3,5,B,?]**

**B → < A_0 B_1 , B_1 A_0 >**

---

**[2,5,B,?*b]**

[2,4] [4,5]     [2,3] [3,5]

A          B

A      B

accept rule(A,B,1)?

y          n

accept item(2,3,A,1)?

y          n

accept item(3,5,B,1)?

y          n

accept item(3,5,B,2)?

y          n

**We want:**

[2,3,A,a*b]

[3,5,B,b*a]

B → < A$_0$ B$_1$ , B$_1$ A$_0$ >

─────────────

[2,5,B,b*b]

**So far we have:**

[2,5,A,a*b]

[3,5,B,b*a]

B → < A$_0$ B$_1$ , B$_1$ A$_0$ >

─────────────

[2,5,B,b*b]

LANGUAGE WEAVER

**[2,4] [4,5]**      **[2,3] [3,5]**

**A**      **B**

**A**      **B**

accept rule(A,B,1)?

**y**      **n**

accept item(2,3,A,1)?

**y**      **n**

accept item(3,5,B,1)?

**y**      **n**

accept item(3,5,B,2)?

**y**      **n**

**this is a search space**

**The Item Creation Problem, rephrased: find the n lowest-cost goal nodes of this search space**

LANGUAGE WEAVER

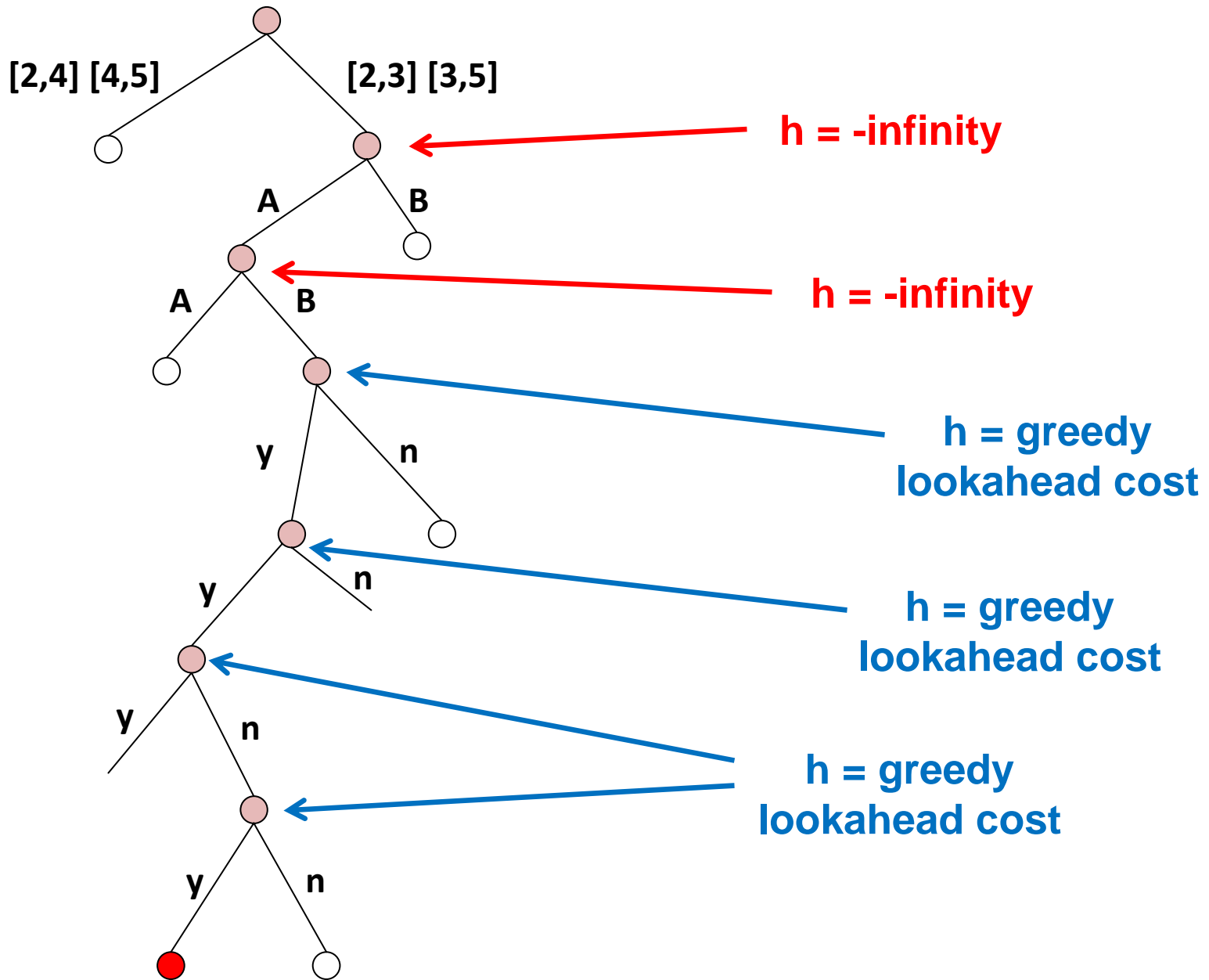[2,4] [4,5]          [2,3] [3,5]

A          B

A          B

y          n

y          n

y          n

y          n
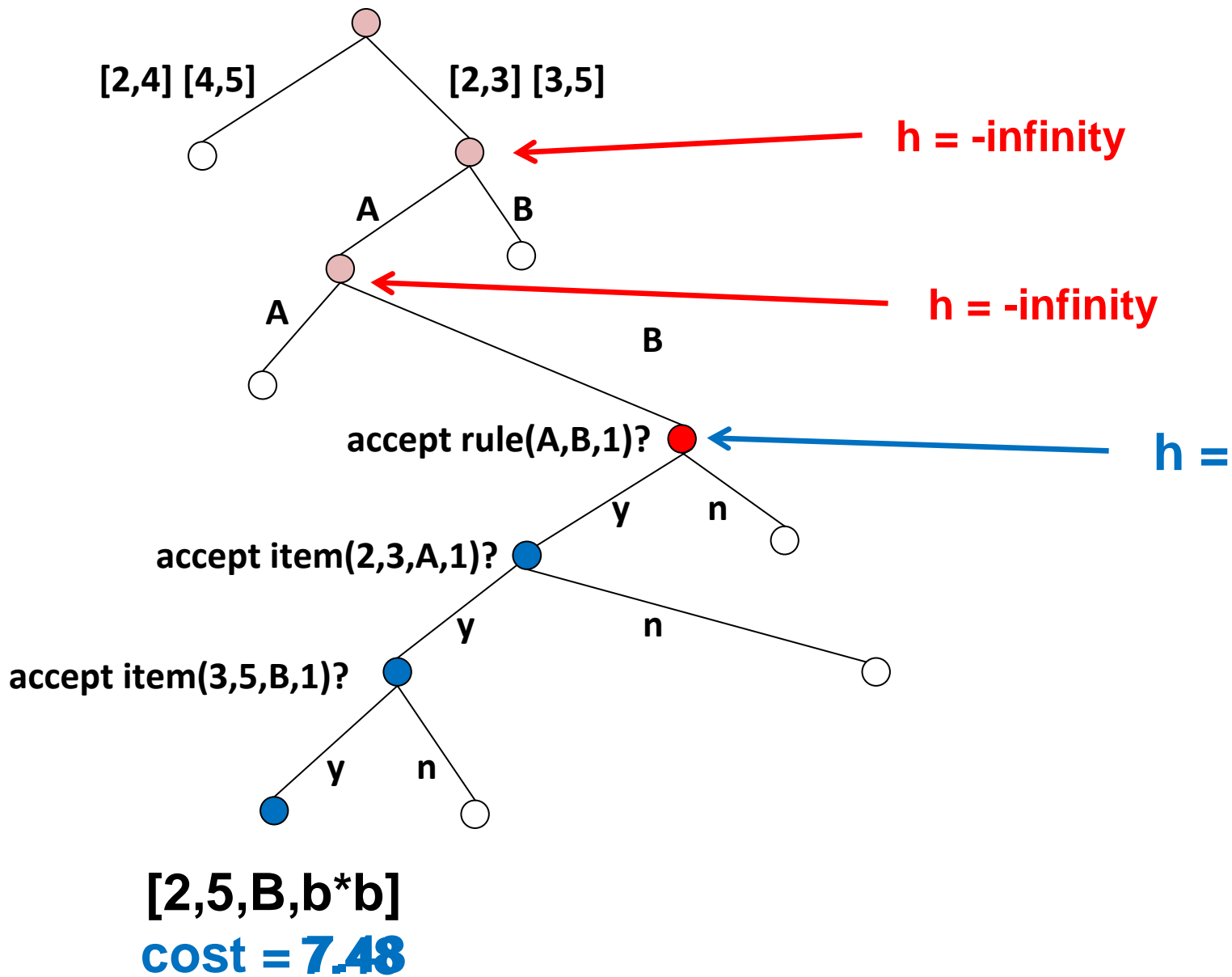
so if we can come up
with lower-bounds on
the best-cost
reachable goal node
from here…
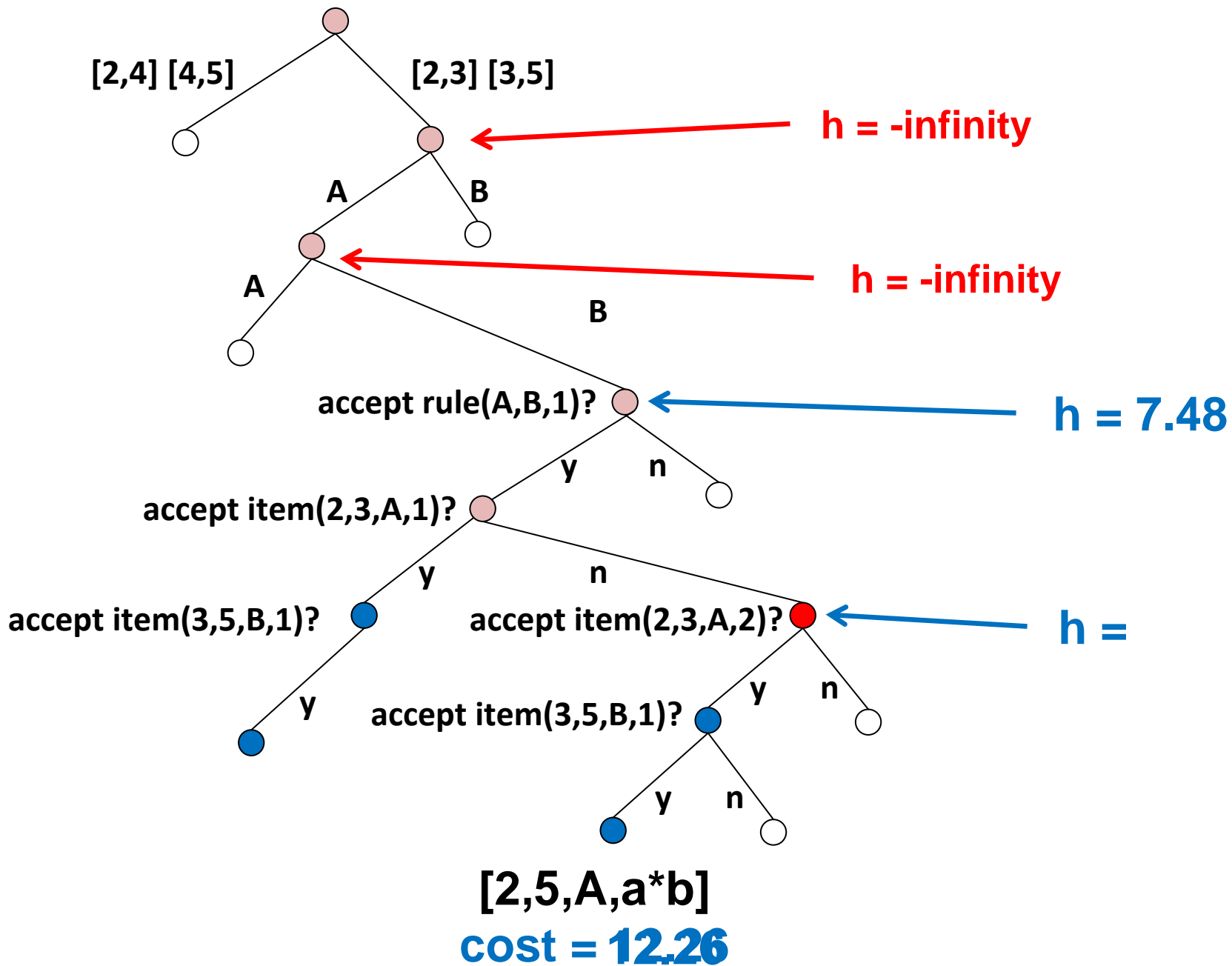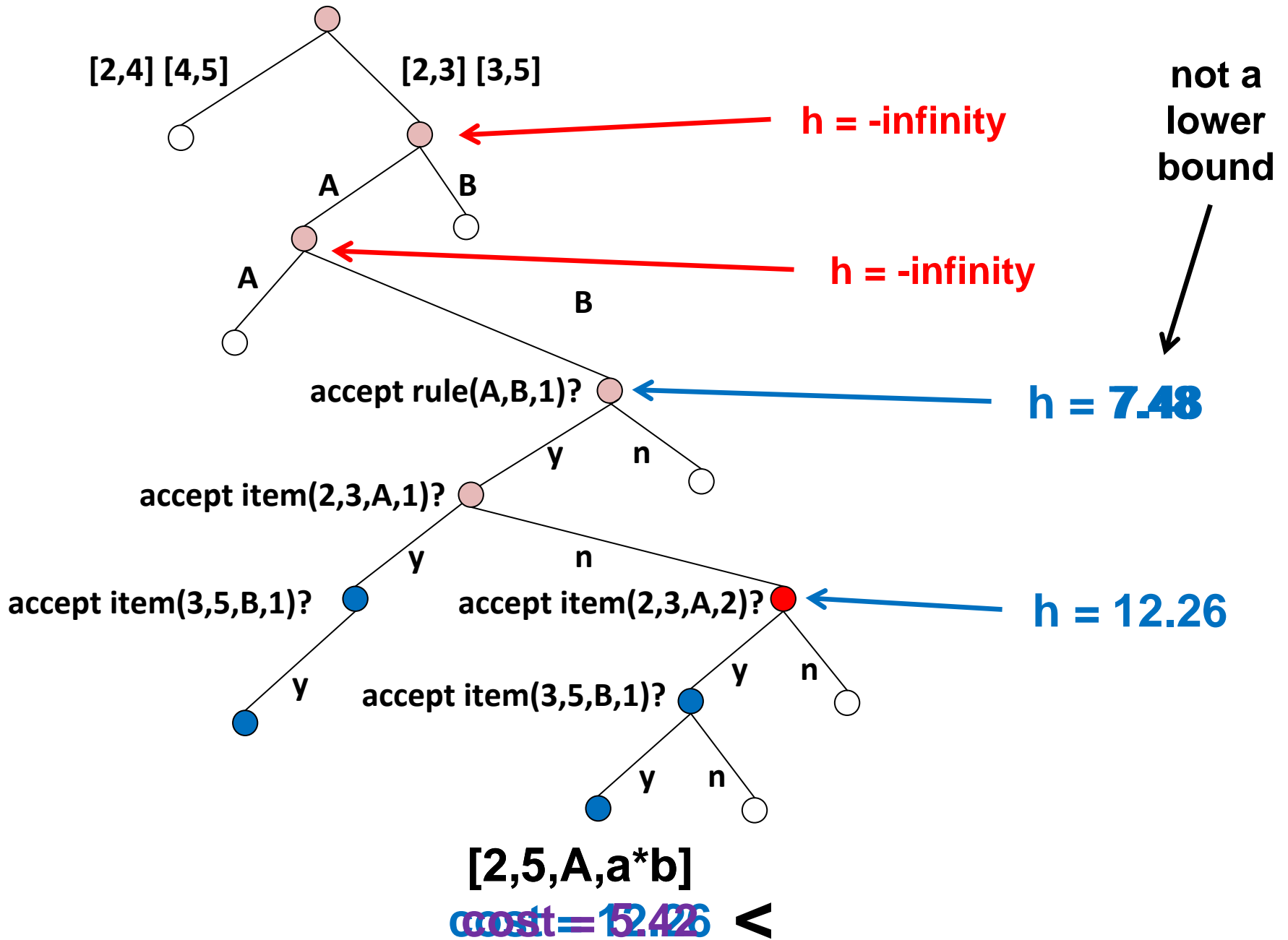…and here

…and here

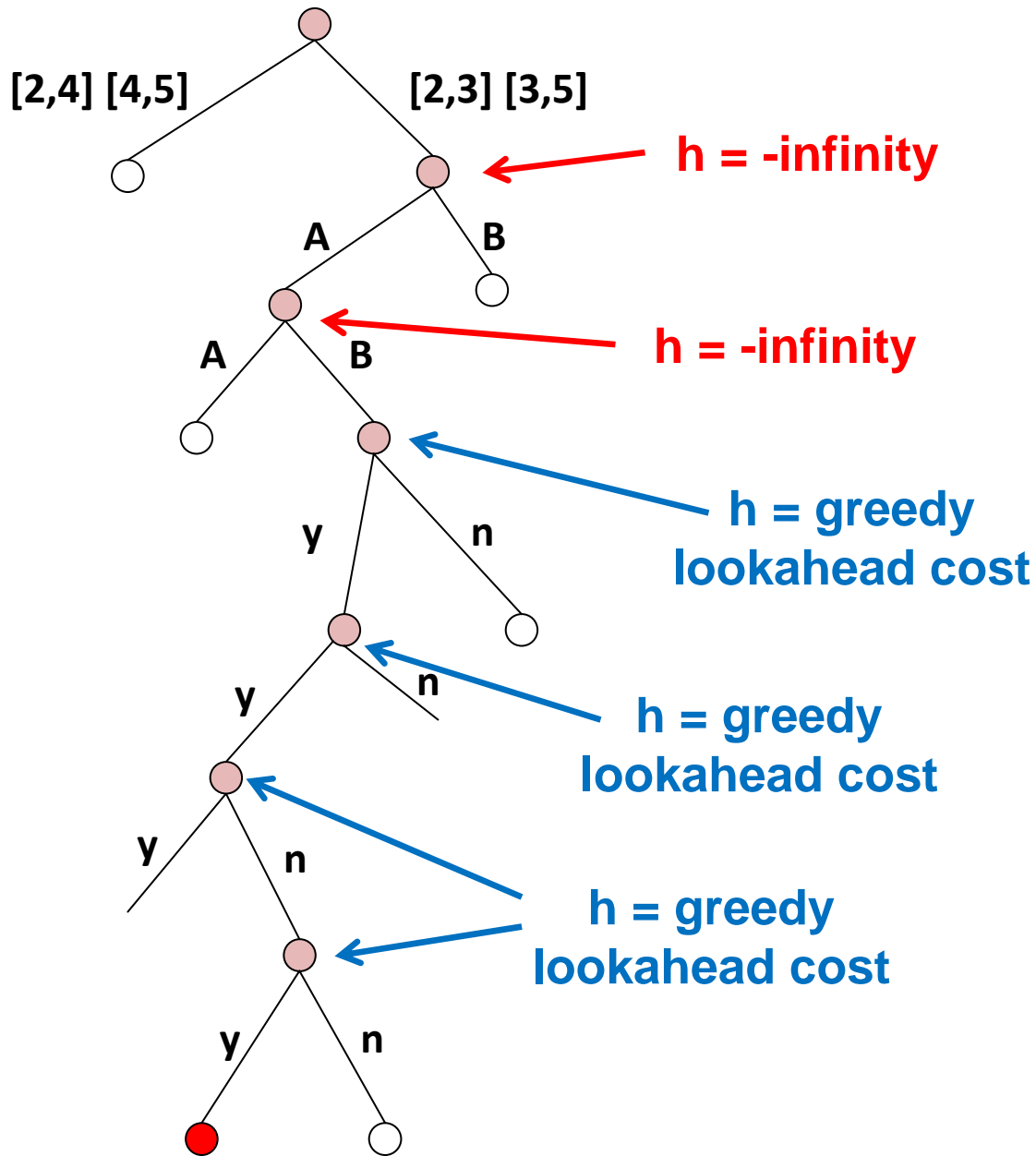…then we can just run
A* on this search
space to find the n
goal nodes of lowest
cost (without
searching the entire
space)

LANGUAGE
WEAVER

[2,4] [4,5]          [2,3] [3,5]

h = -infinity

A          B

A          B

h = -infinity

y          n

h = greedy
lookahead cost

y          n

h = greedy
lookahead cost

y          n

h = greedy
lookahead cost

y          n

LANGUAGE
WEAVER

[2,4] [4,5]   [2,3] [3,5]

h = -infinity

A   B

h = -infinity

A

B

accept rule(A,B,1)?

h =

y   n

accept item(2,3,A,1)?

y   n

accept item(3,5,B,1)?

y   n

[2,5,B,b*b]
cost = 7.48

LANGUAGE WEAVER

[2,4] [4,5]    [2,3] [3,5]

h = -infinity

A    B

A

B

h = -infinity

accept rule(A,B,1)?    h = 7.48

y    n

accept item(2,3,A,1)?

y    n

accept item(3,5,B,1)?    accept item(2,3,A,2)?    h =

y    y    n

accept item(3,5,B,1)?

y    n

[2,5,A,a*b]
cost = 12.26

[2,4] [4,5]   [2,3] [3,5]

h = -infinity

**not a lower bound**

A    B

A

B

accept rule(A,B,1)?    h = 7.48

y    n

accept item(2,3,A,1)?

y    n

accept item(3,5,B,1)?    accept item(2,3,A,2)?    h = 12.26

y    y    n

accept item(3,5,B,1)?

y    n

[2,5,A,a*b]
cost = 15.42 <

[2,4] [4,5]　　　[2,3] [3,5]

h = -infinity

admissible

A　　B

A　　B

h = -infinity

therefore A* will not find the n best solutions, it will only find n good solutions

y　　n

h = greedy lookahead cost

y　　n

h = greedy lookahead cost

y　　n

not admissible

h = greedy lookahead cost

y　　n

# A* search

**VS.**

# Cube pruning

LANGUAGE WEAVER

# Cube pruning begins by forming cubes

[2,3] A
[3,5] A

[2,4] A
[4,5] A

[2,3] B
[3,5] B

[2,4] B
[4,5] B

**and preconditions**

LANGUAGE WEAVER

A*
search

Cube
pruning

**A\* search visits nodes in increasing order of heuristic value**

h = -inf

[2,4] [4,5]          [2,3] [3,5]

h = -inf          h = -inf

A          B          A          B

-inf          -inf          -inf          -inf

A          B          A          B          A          B          A          B

**therefore, it will begin by visiting all nodes with –inf heuristics**

A\*
search

**Cube pruning begins by forming cubes**

| [2,3] A [3,5] B | [2,3] A [3,5] A | [2,4] A [4,5] B | [2,4] A [4,5] A |

| [2,3] B [3,5] A | [2,3] B [3,5] B | [2,4] B [4,5] A | [2,4] B [4,5] B |

Cube
pruning

# A* search visits nodes in order of increasing heuristic value

[2,4] [4,5]        [2,3] [3,5]

A        B        A        B

A    B    A    B    A    B    A    B

# Cube pruning begins by forming cubes

| [2,3] A [3,5] B | [2,3] A [3,5] A | [2,4] A [4,5] B | [2,4] A [4,5] A |

| [2,3] B [3,5] A | [2,3] B [3,5] B | [2,4] B [4,5] A | [2,4] B [4,5] B |

**A* search**

**Cube pruning**

# A* search visits nodes in order of increasing heuristic value

[2,4] [4,5]                                    [2,3] [3,5]

```
                          o
                   /             \
                  o               o
               A /   \ B       A /   \ B
                o     o         o     o
            A /  \ B A / \ B  A / \ B A / \ B
             O    O  O   O    O   O   O    O
```

A
B   A   B   A   B   A   B

## What is a cube?

[2,3] A
[3,5] B

a cube is a set of three axes

| A* search |

| Cube pruning |

LANGUAGE WEAVER

# A* search visits nodes in order of increasing heuristic value

[2,4] [4,5]          [2,3] [3,5]

A          B          A          B

A    B   A    B    A    B   A    B

**A\*
search**

# What is a cube?

sorted by increasing cost

item(2,3,A,3)
item(2,3,A,2)
item(2,3,A,1)

rule(A,B,1)  rule(A,B,2)  rule(A,B,3)  rule(A,B,4)

[2,3]A
[3,5]B

item(3,5,B,1)
item(3,5,B,2)
item(3,5,B,3)

sorted by increasing cost

[2,3] A
[3,5] B

**Cube
pruning**

LANGUAGE
WEAVER

# A* search visits nodes in order of increasing heuristic value



[2,4] [4,5]   [2,3] [3,5]

A   B   A   B

A   B   A   B   A   B   A   B

A* search

# Thus each choice of object from



here  and here

item(2,3,A,3)
item(2,3,A,2)
item(2,3,A,1)

rule(A,B,1)  rule(A,B,2)  rule(A,B,3)  rule(A,B,4)

[2,3]A
[3,5]B

item(3,5,B,1)
item(3,5,B,2)
item(3,5,B,3)

item(2,3,A,2)

item(3,5,B,1)

rule(A,B,4)

here

Cube pruning

LANGUAGE WEAVER

# A* search visits nodes in order of increasing heuristic value

[2,4] [4,5]          [2,3] [3,5]

A          B          A          B

A     B   A     B   A     B   A     B

## …creates a new item for [2,5]

item(2,3,A,3)
item(2,3,A,2)
item(2,3,A,1)

rule(A,B,1)  rule(A,B,2)  rule(A,B,3)  rule(A,B,4)

[2,3]A
[3,5]B

item(3,5,B,1)
item(3,5,B,2)
item(3,5,B,3)

[2,3,A,a*2]

[3,5,B,b]

$B \rightarrow$ < $A_0(B_1,B_1)$ $A_0$ >

$[2,5,B,b*b]$

A*
search

Cube
pruning

LANGUAGE WEAVER

# A* search visits nodes in order of increasing heuristic value

[2,4] [4,5]          [2,3] [3,5]

A          B          A          B

A     B   A     B    A     B    A     B

**If we take the best representative from each axis (i.e. the "1-1-1")…**

item(2,3,A,3)
item(2,3,A,2)
item(2,3,A,1)

rule(A,B,1)  rule(A,B,2)  rule(A,B,3)  rule(A,B,4)

[2,3]A
[3,5]B

item(3,5,B,1)
item(3,5,B,2)
item(3,5,B,3)

cost( new item ) =
    …then we expect
    cost( subitem1 )
    the resulting item
+ cost( subitem2 )
    to have a low cost,
+ cost( rule )
        since:
+ interaction( rule, subitems )

A*
search

Cube
pruning

LANGUAGE WEAVER

**A\* search visits nodes in order of increasing heuristic value**

[2,4] [4,5]          [2,3] [3,5]

A          B          A          B

A     B    A     B    A     B    A     B

**A\* search**

**Though we are not guaranteed this, because:**

*item(2,3,A,3)*
*item(2,3,A,2)*
*item(2,3,A,1)*

*rule(A,B,1)* *rule(A,B,2)* *rule(A,B,3)* *rule(A,B,4)*

[2,3]A
[3,5]B

*item(3,5,B,1)*
*item(3,5,B,2)*
*item(3,5,B,3)*

cost( **new item** ) =
 cost( **subitem1** )
+ cost( **subitem2** )
+ cost( **rule** )
+ interaction( rule, subitems )

**this cost is not monotonic**

**Cube pruning**

LANGUAGE WEAVER

# A* search visits nodes in order of increasing heuristic value

[2,4] [4,5]          [2,3] [3,5]

A          B          A          B

A          B    A    B    A    B    A    B

# Cube pruning proceeds by creating the 1-1-1 item of every cube

| 111 | 111 | 111 | 111 |
|-----|-----|-----|-----|
| 2.4 | 4.6 | 7.9 | 3.2 |

| 111 | 111 | 111 | 111 |
|-----|-----|-----|-----|
| 5.5 | 9.2 | 6.2 | 4.4 |

## and scoring them

A* search

Cube pruning

LANGUAGE WEAVER

# Meanwhile, A* search has scored its frontier nodes

[2,4] [4,5]          [2,3] [3,5]

A          B          A          B

A     B     A     B     A     B     A     B

accept rule(A,B,1)

accept item(2,3,A,1)

accept item(3,5,B,1)

**A* search**

**111**

**2.4**

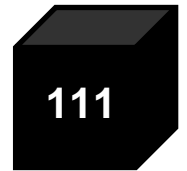# Cube pruning proceeds by creating the 1-1-1 item of every cube

**111**    **111**    **111**    **111**

**2.4**    **4.6**    **7.9**    **3.2**

**111**    **111**    **111**    **111**

**5.5**    **9.2**    **6.2**    **4.4**

## and scoring them

**Cube pruning**

LANGUAGE WEAVER

# Meanwhile, A* search has scored its frontier nodes

[2,4] [4,5]     [2,3] [3,5]

A     B     A     B

A   B   A   B   A   B   A   B

**3.2   7.9   6.2   4.4   4.6   2.4   5.5   9.2**

A*
search

# Cube pruning proceeds by creating the 1-1-1 item of every cube

| 111 | 111 | 111 | 111 |
|-----|-----|-----|-----|
| **2.4** | **4.6** | **7.9** | **3.2** |

| 111 | 111 | 111 | 111 |
|-----|-----|-----|-----|
| **5.5** | **9.2** | **6.2** | **4.4** |

## and scoring them

Cube
pruning

LANGUAGE
WEAVER

**Meanwhile, A\* search has scored its frontier nodes**

[2,4] [4,5]　　　　[2,3] [3,5]

A　　B　　　A　　B

A　　B　A　　B　A　　B　A　　B

3.2　7.9　6.2　4.4　4.6　2.4　5.5　9.2

| A\* |
|-----|
| **A\*** |
| **search** |

**At this point, And cube generates its "best off" items. as created.**

| 121 | 112 |
|-----|-----|
| 7.1 | 3.6 |

| 211 | 111 | 111 | 111 |
|-----|-----|-----|-----|
| 2.8 | 4.6 | 7.9 | 3.2 |

| 111 | 111 | 111 | 111 |
|-----|-----|-----|-----|
| 5.5 | 9.2 | 6.2 | 4.4 |

**It keeps this item.**

**KEPT**

| **Cube** |
|----------|
| **pruning** |

# A* search continues to visit nodes in increasing order of heuristic value

[2,4] [4,5]          [2,3] [3,5]

A          B          A          B

A          B     A          B     A          B     A          B

3.2     7.9     6.2     4.4     4.6     2.4     5.5     9.2

**A* search**

# And generates its "one-off" items.

| | |
|---|---|
| **121** | **112** |
| 7.1 | 3.6 |

| | | | |
|---|---|---|---|
| **211** | **111** | **111** | **111** |
| 4.8 | 4.6 | 7.9 | 3.2 |

| | | | |
|---|---|---|---|
| **111** | **111** | **111** | **111** |
| 5.5 | 9.2 | 6.2 | 4.4 |

## It keeps this item.

**KEPT** 111

**Cube pruning**

LANGUAGE WEAVER

# A* search continues to visit nodes in increasing order of heuristic value

[2,4] [4,5]      [2,3] [3,5]

A    B      A    B

A   B   A   B    A   B   A   B

3.2   7.9   6.2   4.4   4.6      5.5    9.2

y   n

2.4   4.8
111   211

**A* search**

# And generates its "one-off" items.

| 121 | 112 |
|-----|-----|
| 7.1 | 3.6 |

| 211 | 111 | 111 | 111 |
|-----|-----|-----|-----|
| 4.8 | 4.6 | 7.9 | 3.2 |

| 111 | 111 | 111 | 111 |
|-----|-----|-----|-----|
| 5.5 | 9.2 | 6.2 | 4.4 |

## It keeps this item.

KEPT

**Cube pruning**

**A\* search continues to visit nodes in increasing order of heuristic value**

[2,4] [4,5]  [2,3] [3,5]

A  B  A  B

A  B  A  B  A  B  A  B

3.2  7.9  6.2  4.4  4.6  5.5  9.2

y  n

2.4  4.8

**A\* search**

**And generates its "one-off" items.**

| 121 | 112 |
| 7.1 | 3.6 |

| 211 | 111 | 111 | 111 |
| 4.8 | 4.6 | 7.9 | 3.2 |

| 111 | 111 | 111 | 111 |
| 5.5 | 9.2 | 6.2 | 4.4 |

**It keeps this item.**

KEPT

**Cube pruning**

LANGUAGE WEAVER

**A\* search continues to visit nodes in increasing order of heuristic value**

[2,4] [4,5]          [2,3] [3,5]

A          B          A          B

A     B    A     B    A     B    A     B

3.2   7.9  6.2   4.4  4.6       5.5   9.2

y          n

4.8

y          n

111  2.4        7.1  121

**A\* search**

**And generates its "one-off" items.**

121          112
7.1          3.6

211          111          111          111
4.8          4.6          7.9          3.2

111          111          111          111
5.5          9.2          6.2          4.4

**It keeps this item.**

KEPT  111

**Cube pruning**

LANGUAGE WEAVER

**A\* search continues to visit nodes in increasing order of heuristic value**

[2,4] [4,5]          [2,3] [3,5]

A          B          A          B

A     B     A     B     A     B     A     B

3.2    7.9    6.2    4.4    4.6         5.5    9.2

y          n

4.8

y          n

2.4          7.1

**A\* search**

---



**And generates its "one-off" items.**

| 121 | 112 |
|---|---|
| 7.1 | 3.6 |

| 211 | 111 | 111 | 111 |
|---|---|---|---|
| 4.8 | 4.6 | 7.9 | 3.2 |

| 111 | 111 | 111 | 111 |
|---|---|---|---|
| 5.5 | 9.2 | 6.2 | 4.4 |

**It keeps this item.**

KEPT

**Cube pruning**

LANGUAGE WEAVER

# A* search continues to visit nodes in increasing order of heuristic value

[2,4] [4,5]          [2,3] [3,5]

A          B          A          B

A     B     A     B     A     B     A     B

3.2    7.9    6.2    4.4    4.6          5.5    9.2

y          n
                    4.8

y          n
                    7.1

y          n
12.4               3.6 112

**A* search**

## And generates its "one-off" items.

| 121 | 112 |
|-----|-----|
| 7.1 | 3.6 |

| 211 | 111 | 111 | 111 |
|-----|-----|-----|-----|
| 4.8 | 4.6 | 7.9 | 3.2 |

| 111 | 111 | 111 | 111 |
|-----|-----|-----|-----|
| 5.5 | 9.2 | 6.2 | 4.4 |

## It keeps this item.

111 KEPT

**Cube pruning**

LANGUAGE WEAVER

# A* search continues to visit nodes in increasing order of heuristic value

[2,4] [4,5]        [2,3] [3,5]

A        B        A        B

A    B    A    B    A    B    A    B

3.2   7.9   6.2   4.4   4.6        5.5        9.2

y        n
                    4.8

y        n
                    7.1

this is a goal node →

y        n
2.4        3.6

A* search

# And generates its "one-off" items.

| 121 | 112 |
| --- | --- |
| 7.1 | 3.6 |

| 211 | 111 | 111 | 111 |
| --- | --- | --- | --- |
| 4.8 | 4.6 | 7.9 | 3.2 |

| 111 | 111 | 111 | 111 |
| --- | --- | --- | --- |
| 5.5 | 9.2 | 6.2 | 4.4 |

## It keeps this item.

KEPT

Cube pruning

LANGUAGE WEAVER

**A\* search continues to visit nodes in increasing order of heuristic value**

[2,4] [4,5]     [2,3] [3,5]

A    B     A    B

A    B    A    B    A    B    A    B

3.2   7.9   6.2   4.4   4.6   5.5   9.2

y    n
4.8

y    n
7.1

n
3.6

**So A\* keeps this item.**

this is a goal this item.

**KEPT** 121

**A\* search**

---

**And generates its "one-off" items.**

121     112
7.1     3.6

211     111     111     111
4.8     4.6     7.9     3.2

111     111     111     111
5.5     9.2     6.2     4.4

**It keeps this item.**

**KEPT** 111

**Cube pruning**

LANGUAGE WEAVER

**A\* search**

**(+ node tying)**

**=**

**Cube pruning**

LANGUAGE WEAVER

[2,4] [4,5]    [2,3] [3,5]

h = -infinity

?

A    B

h = -infinity

A    B

accept rule(A,B,1)?

h = greedy
lookahead cost

y    n

accept item(2,3,A,1)?

Cube pruning was
specifically designed for
hierarchical phrase-based
MT

y    n

h = greedy
lookahead cost

accept item(3,5,B,1)?

But say our use case was
string-to-tree MT, in the style
of (Galley et al 2006)

y    n

which has a small
number of distinct
postconditions

accept item(3,5,B,2)?

h = greedy
lookahead cost

y    n

# Average number of search nodes visited, per sentence

## Arabic-English NIST 2008

| Nodes by Type | Cube Pruning |
|---|---|
| subspan | 12936 |
| precondition | **851458** |
| rule | 33734 |
| item | 119703 |
| goal | 74618 |
| TOTAL | **1092449** |
| BLEU | **38.33** |

**the early nodes with infinite heuristics dominate the search time**

LANGUAGE WEAVER

[2,4] [4,5]          [2,3] [3,5]

h = something better

A          B

A          B

accept rule(A,B,1)?

h = something better

y          n

accept item(2,3,A,1)?

h = greedy lookahead cost

y          n

accept item(3,5,B,1)?

h = greedy lookahead cost

y          n

accept item(3,5,B,2)?

h = greedy lookahead cost

y          n

LANGUAGE WEAVER

# Number of search nodes visited

## Arabic-English NIST 2008

| Nodes by Type | Cube Pruning | Augmented CP |
|---|---|---|
| subspan | 12936 | 12792 |
| precondition | **851458** | **379954** |
| rule | 33734 | 33331 |
| item | 119703 | 118889 |
| goal | 74618 | 74159 |
| TOTAL | **1092449** | **619125** |
| BLEU | **38.33** | **38.22** |

LANGUAGE WEAVER

# Tradeoff curves

## Arabic-English NIST 2008

[2,4] [4,5]          [2,3] [3,5]

h = admissible

A          B

A          B

h = admissible

Cube pruning becomes **exact**.

accept rule(A,B,1)?

We found that our exact version of cube pruning had a somewhat better quality curve than the original inexact version of cube pruning.

h = admissible
lookahead cost

y          n

accept item(2,3,A,1)?

This is full reepening but a similar/predictable result to the

h = guessible
lookahead cost

y          n

accept item(3,5,B,1)?

y          n

However it was not as effective as our "augmented" version of cube pruning.

h = guessible
lookahead cost

accept item(3,5,B,2)?

h = guessible
lookahead cost

y          n

LANGUAGE WEAVER

# What **L**esson

**does cube pruning teach us?**

LANGUAGE
WEAVER

[2,4] [4,5]   [2,3] [3,5]

A        B

A      B

accept rule(A,B,1)?

y        n

accept item(2,3,A,1)?

y

n

accept item(3,5,B,1)?

y      n

accept item(3,5,B,2)?

y      n

It tells us that it is useful to frame the CYK Item Generation Problem as a Deuristic search realization problem.

Once this search realization is made, we suddenly have many more avenues available this is a to us, when search space implementing a CKY decoder for a particular use case

h = something better

h = something better

h = greedy lookahead cost

We can change the heuristics.

h = greedy lookahead cost

h = greedy lookahead cost

LANGUAGE WEAVER

**We can change
the search algorithm.**

LANGUAGE
WEAVER

**For instance, instead of A*…**

**We can change
the search algorithm.**

LANGUAGE
WEAVER

**For instance, instead of A*…**

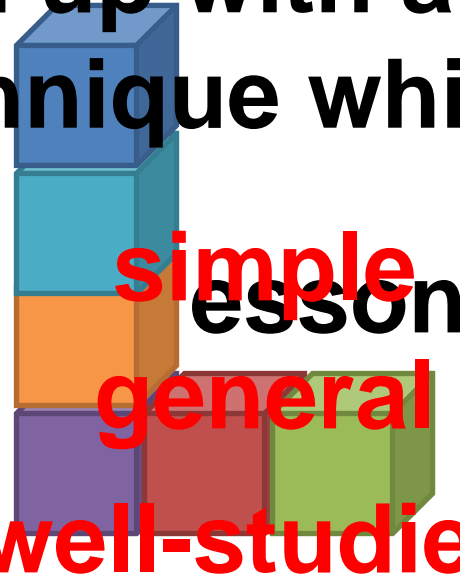**…we could try a depth-first strategy like depth-first branch-and-bound, and take advantage of its anytime properties.**

**We can change the search algorithm.**

LANGUAGE
WEAVER

We can change
the search algorithm.
the search space.

LANGUAGE
WEAVER

# We end up with a speedup technique which is:

**What** **Lesson**

**simple**

**general**

**well-studied**

**easily adaptable to**

**new use cases**

# does cube pruning teach us?

# Thank you.

## Questions?

LANGUAGE
WEAVER